

## 0.1 Conditional preference-based energy learning with HTE embeddings and an LP economic oracle

**Scenario encoding with a Hierarchical Temporal Encoder (HTE).** Each scenario  $x$  (multilayer grid graph, exogenous drivers, and time series) is first encoded by a *Hierarchical Temporal Encoder* (HTE) into a context representation

$$h = \text{HTE}_\phi(x), \quad h \in \mathbb{R}^d, \quad (1)$$

where  $\phi$  denotes the HTE parameters (trained beforehand or jointly, depending on the experimental setting). In practice,  $h$  may include multi-scale components (e.g., zone- and system-level embeddings) and temporal summaries; throughout this section we use  $h$  as a generic shorthand for the HTE-provided conditioning signal.

**Decision space.** The operational decision is represented by a high-dimensional binary vector  $u \in \{0, 1\}^M$  collecting unit-commitment (UC) flags, discrete demand-response (DR) activations, and discrete storage mode indicators across space and time. For sampling, we often operate in a continuous relaxation  $\tilde{u} \in (0, 1)^M$  (or in logit space) and binarize only before the physics-based refinement stage.

**Conditional energy-based model (EBM).** We learn a *conditional* energy function

$$E_\theta(u | h) \in \mathbb{R}, \quad (2)$$

where  $\theta$  are the EBM parameters and  $h$  is the HTE embedding of the scenario. The EBM defines an implicit conditional distribution over discrete decisions:

$$p_\theta(u | h) \propto \exp(-E_\theta(u | h)). \quad (3)$$

Low-energy configurations are intended to correspond to low-cost operational strategies for scenario context  $h$ .

**Implicit policy via normalized Langevin sampling.** Given  $h$ , we generate a set of  $K$  candidate configurations by running a stochastic sampler targeting (3). Let  $\mathcal{S}_\theta$  denote a normalized Langevin procedure operating in a relaxed space:

$$\{\tilde{u}^{(k)}\}_{k=1}^K \sim \mathcal{S}_\theta(h), \quad \tilde{u}^{(k)} \in (0, 1)^M, \quad (4)$$

and convert relaxed samples into binary decisions

$$u^{(k)} = \text{Bin}(\tilde{u}^{(k)}) \in \{0, 1\}^M, \quad (5)$$

where  $\text{Bin}(\cdot)$  is thresholding or Bernoulli sampling.

**Hierarchical feasibility decoder and LP worker (economic oracle).** Each discrete candidate  $u^{(k)}$  is passed through a hierarchical decoder  $\mathcal{D}$  that produces a coherent continuous initialization (dispatch, storage trajectories, DR profiles, and flows) consistent with  $u^{(k)}$ :

$$y_{\text{dec}}^{(k)} = \mathcal{D}(x, u^{(k)}). \quad (6)$$

Then, an LP worker  $\mathcal{W}$  *hard-fixes* the discrete decisions  $u^{(k)}$  and solves a continuous dispatch problem to validate feasibility and compute the realized operational cost:

$$(\text{feas}^{(k)}, y_*^{(k)}, C^{(k)}) = \mathcal{W}(x, u^{(k)}, y_{\text{dec}}^{(k)}), \quad (7)$$

where  $\text{feas}^{(k)} \in \{0, 1\}$  denotes feasibility and  $C^{(k)} \in \mathbb{R}_+$  is the resulting cost (including penalty terms such as VOLL when applicable). Crucially,  $\mathcal{W}$  is treated as a *non-differentiable* physics-aware oracle: it provides exact economic feedback, but gradients are *not* propagated through the LP.

**Preference signal induced by realized costs.** Among feasible candidates we select the lowest-cost solution

$$k^* \in \arg \min_{k: \text{feas}^{(k)}=1} C^{(k)}, \quad \hat{u} := u^{(k^*)}. \quad (8)$$

When available, we also obtain a reference solution  $u^+$  from the MILP oracle (optimal or best incumbent under a time limit) with cost  $C^+$ . The LP worker thus induces an ordering over decisions:

$$u^{(i)} \succ u^{(j)} \iff C^{(i)} < C^{(j)}. \quad (9)$$

Learning then aims at shaping the energy landscape so that lower-cost decisions receive lower energy under the same conditioning  $h$ .

**Preference-based objective (conditional energy shaping).** We train the conditional energy  $E_\theta(\cdot | h)$  using a margin-ranking loss comparing the MILP reference  $u^+$  to hard candidates produced by the pipeline. Let  $\mathcal{K}$  be a set of “hard negatives” (e.g., feasible candidates with high realized costs, or a diverse subset). We minimize

$$\mathcal{L}_{\text{rank}}(\theta) = \frac{1}{|\mathcal{K}|} \sum_{k \in \mathcal{K}} \max\left(0, m + E_\theta(u^+ | h) - E_\theta(u^{(k)} | h)\right), \quad (10)$$

where  $m > 0$  is a margin. To emphasize rare but catastrophic failures, we optionally use a cost-aware weighting based on the realized cost gap:

$$w_k = \text{clip}\left(\log(1 + (C^{(k)} - C^+)_+), 0, w_{\max}\right), \quad (11)$$

and define

$$\mathcal{L}_{\text{w-rank}}(\theta) = \frac{1}{|\mathcal{K}|} \sum_{k \in \mathcal{K}} (1 + \alpha w_k) \max\left(0, m + E_\theta(u^+ | h) - E_\theta(u^{(k)} | h)\right), \quad (12)$$

with  $\alpha \geq 0$ .

**Gradient flow and training protocol.** The full pipeline (Langevin sampler  $\rightarrow$  decoder  $\rightarrow$  LP worker) is used to generate economically meaningful hard candidates and costs. Gradients are computed *only* through the EBM evaluations  $E_\theta(\cdot | h)$  in Eqs. (10)–(12). In particular, we do not differentiate through the decoder  $\mathcal{D}$ , the LP worker  $\mathcal{W}$ , nor the MILP oracle. The conditioning by HTE embeddings  $h = \text{HTE}_\phi(x)$  ensures that the learned energy landscape adapts to scenario-specific spatio-temporal patterns (demand, VRE, congestion, storage tension), enabling fast inference via sampling while retaining physics-aware economic evaluation through the LP worker.

---

**Algorithm 1:** Conditional preference-based EBM training (HTE-conditioned) with an LP economic oracle

---

**Input:** Scenario dataset  $\{x_n\}$ ; HTE encoder  $\text{HTE}_\phi$ ; EBM  $E_\theta$   
**Input:** Sampler  $\mathcal{S}_\theta$ ; decoder  $\mathcal{D}$ ; LP worker  $\mathcal{W}$   
**Input:** Hyperparameters:  $K$  candidates, margin  $m$ , weighting  $\alpha$

1 **for** each minibatch of scenarios  $x$  **do**

2   Compute conditioning embeddings  $h \leftarrow \text{HTE}_\phi(x)$ ;

3   Sample  $\{\tilde{u}^{(k)}\}_{k=1}^K \sim \mathcal{S}_\theta(h)$ ;

4   Binarize  $u^{(k)} \leftarrow \text{Bin}(\tilde{u}^{(k)})$ ;

5   **for**  $k = 1$  **to**  $K$  **do**

6      $y_{\text{dec}}^{(k)} \leftarrow \mathcal{D}(x, u^{(k)})$ ;

7      $(\text{feas}^{(k)}, C^{(k)}) \leftarrow \mathcal{W}(x, u^{(k)}, y_{\text{dec}}^{(k)})$ ;

8   Select a hard-negative set  $\mathcal{K}$  from feasible candidates (e.g., high-cost feasible);

9   Retrieve MILP reference  $(u^+, C^+)$  for these scenarios when available;

10   Compute  $\mathcal{L}_{\text{w-rank}}$  using energies  $E_\theta(\cdot | h)$  only;

11   Update  $\theta$  by backpropagation through  $E_\theta$  (no gradients through  $\mathcal{D}$  or  $\mathcal{W}$ );

---